



GENERATING REAL-WORLD AND SYNTHETIC GRAPH DATASET FOR GNN APPLICATION

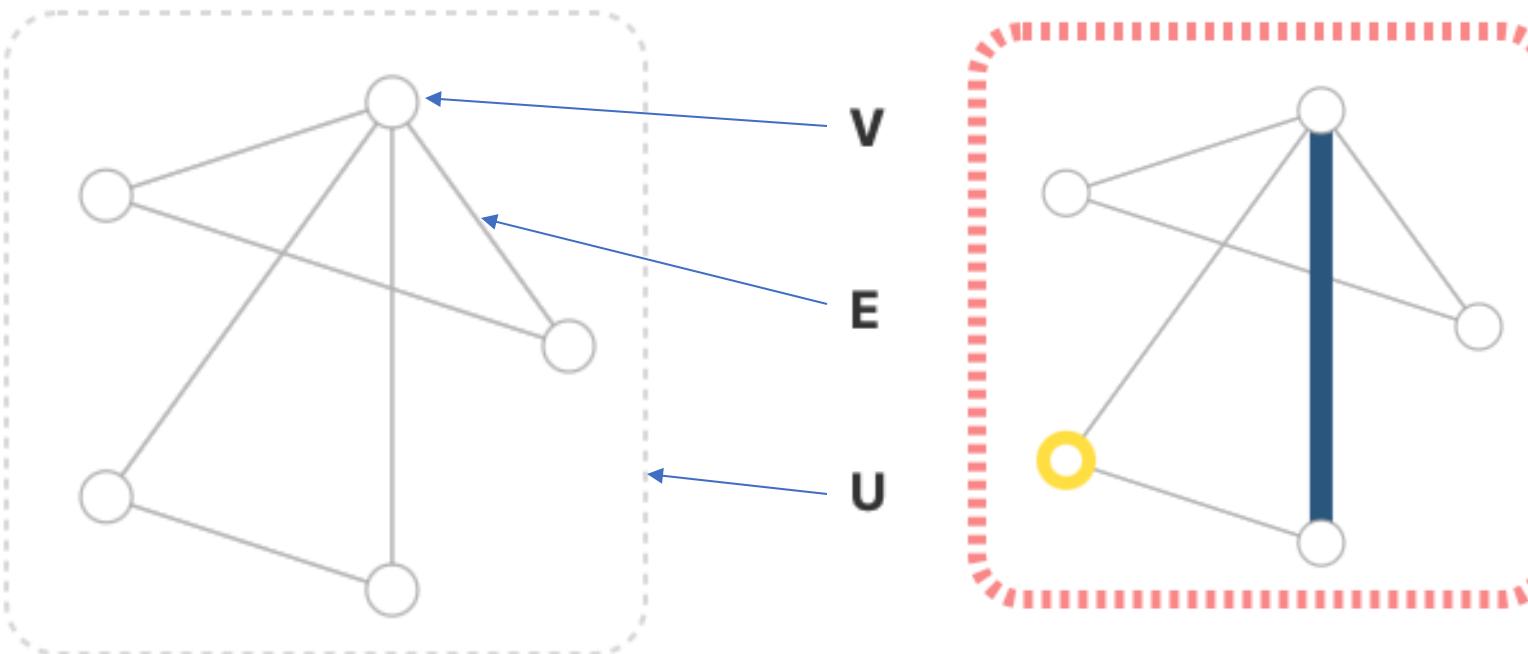
Arpandeep Khatua

Adviser: Professor Wen-Mei Hwu

Mentor: Vikram Sharma Mailthody

20th April 2022

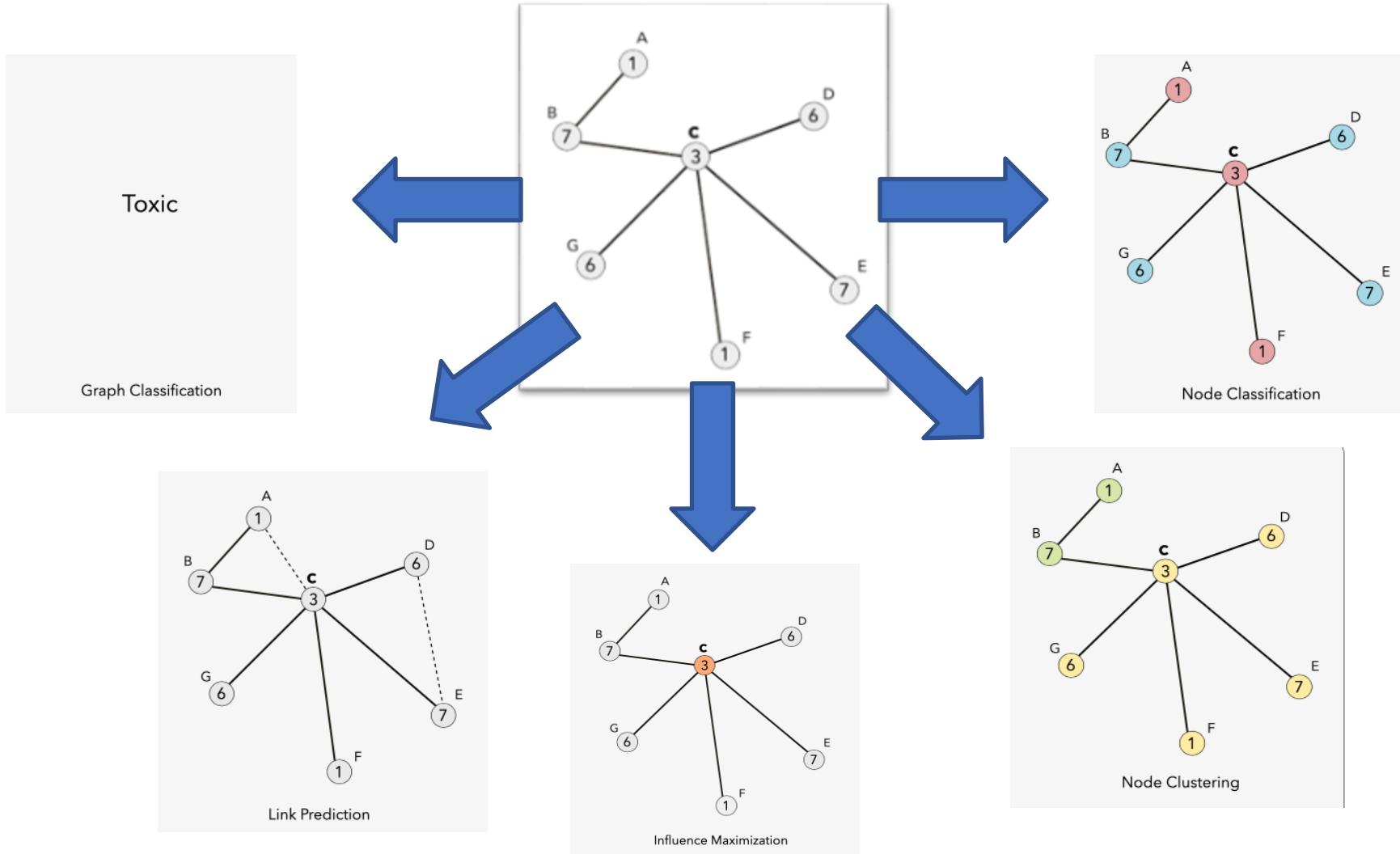
WHAT ARE GRAPHS AND WHY DO WE NEED THEM?



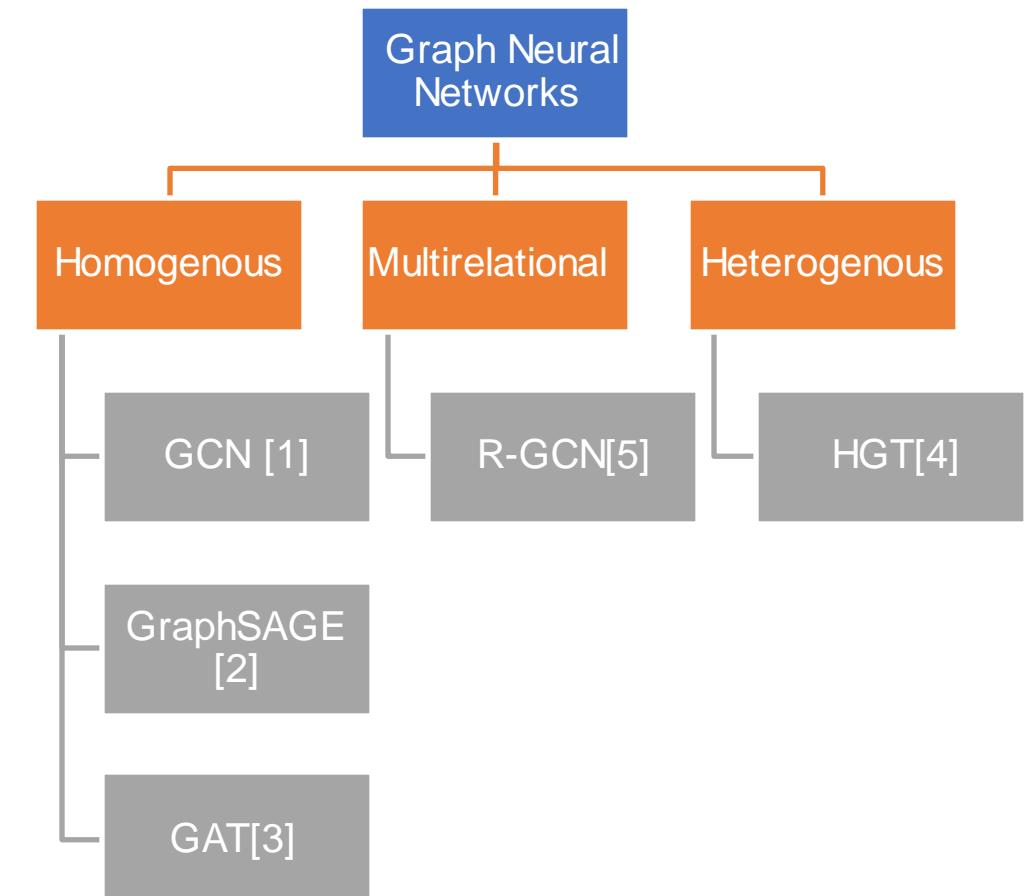
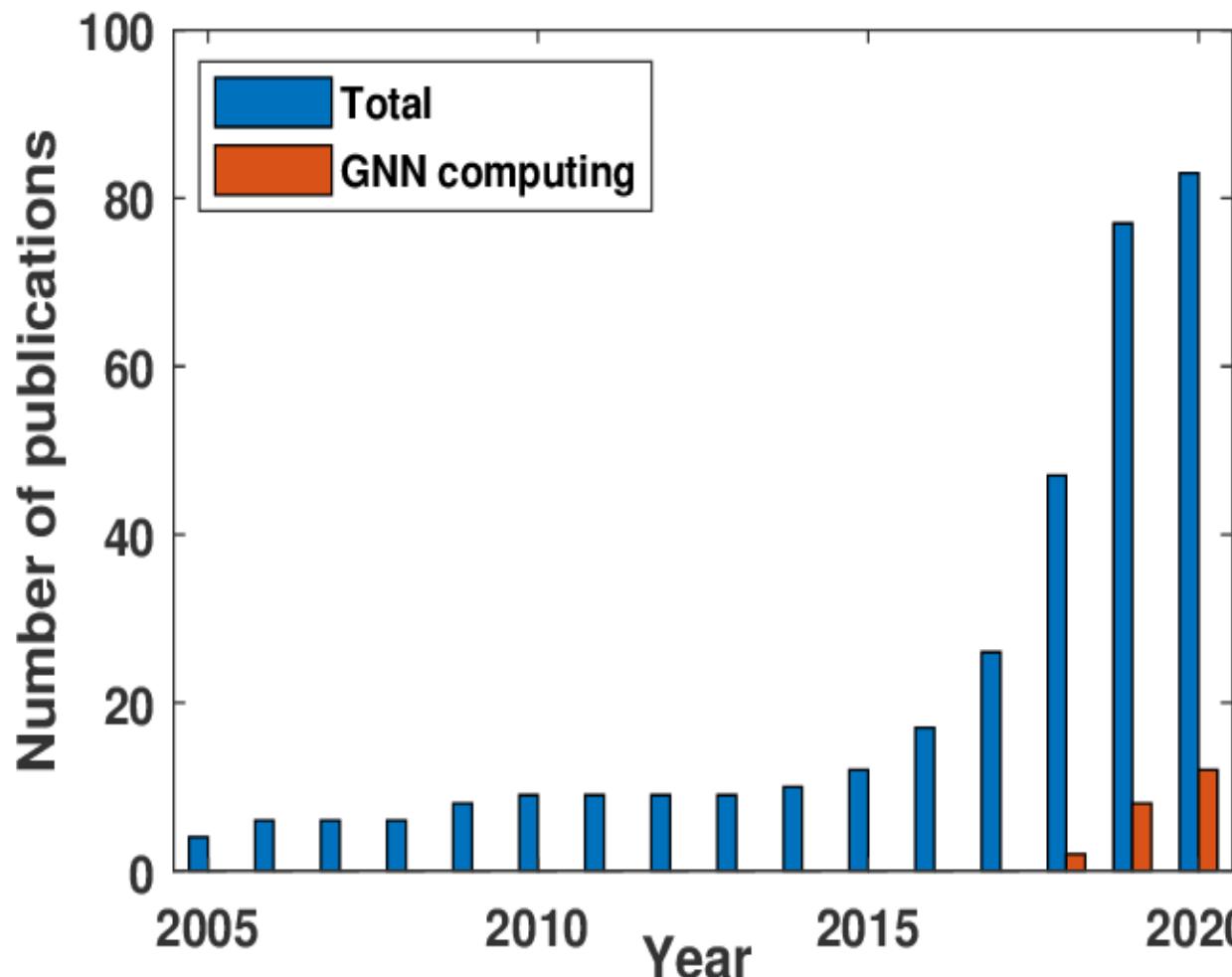
- Vertex (or node) embedding
 - Edge (or link) attributes and embedding
 - Global (or master node) embedding
- Each item is accompanied by a small horizontal bar chart:
- Vertex embedding: A yellow bar chart consisting of several short segments.
 - Edge attributes/embedding: A blue bar chart consisting of several short segments.
 - Global embedding: A pink bar chart consisting of several short segments.

OKAY... SO WHAT CAN WE DO WITH GRAPHS?

There are three general types of prediction tasks on graphs: **graph-level**, **node-level**, and **edge-level**.



DEEP LEARNING IN GRAPHS



MAIN COMPONENTS OF MY RESEARCH WORK

As a part of the thesis, I had two main tasks:

1. To understand the latest GNN models and their limitations.
2. Generating a large graph dataset for both systems designer and large graph datasets.

Our final goal is to create a large graph dataset. However, this is not possible without understanding how GNN Models and previous graph datasets were made and work.

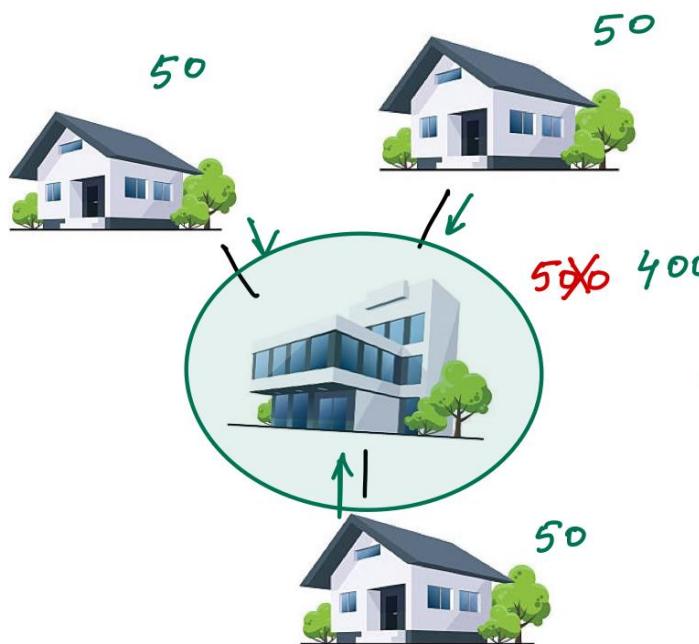
DEEP LEARNING IN GRAPHS



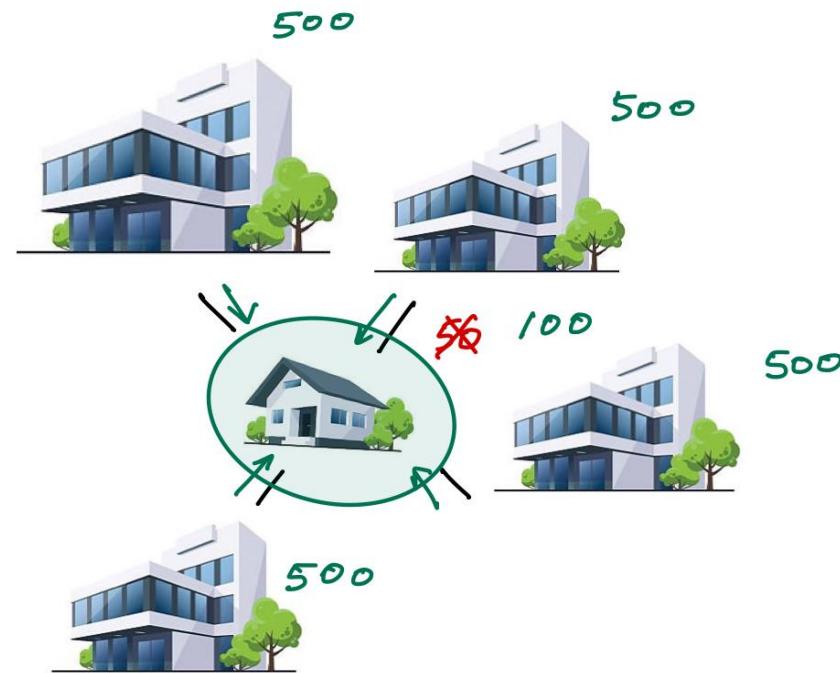
Basic House $\xrightarrow{\text{Emb}}$ e $P(e) = 50$



Mansion $\xrightarrow{\text{Emb}}$ X $P(X) = 500$

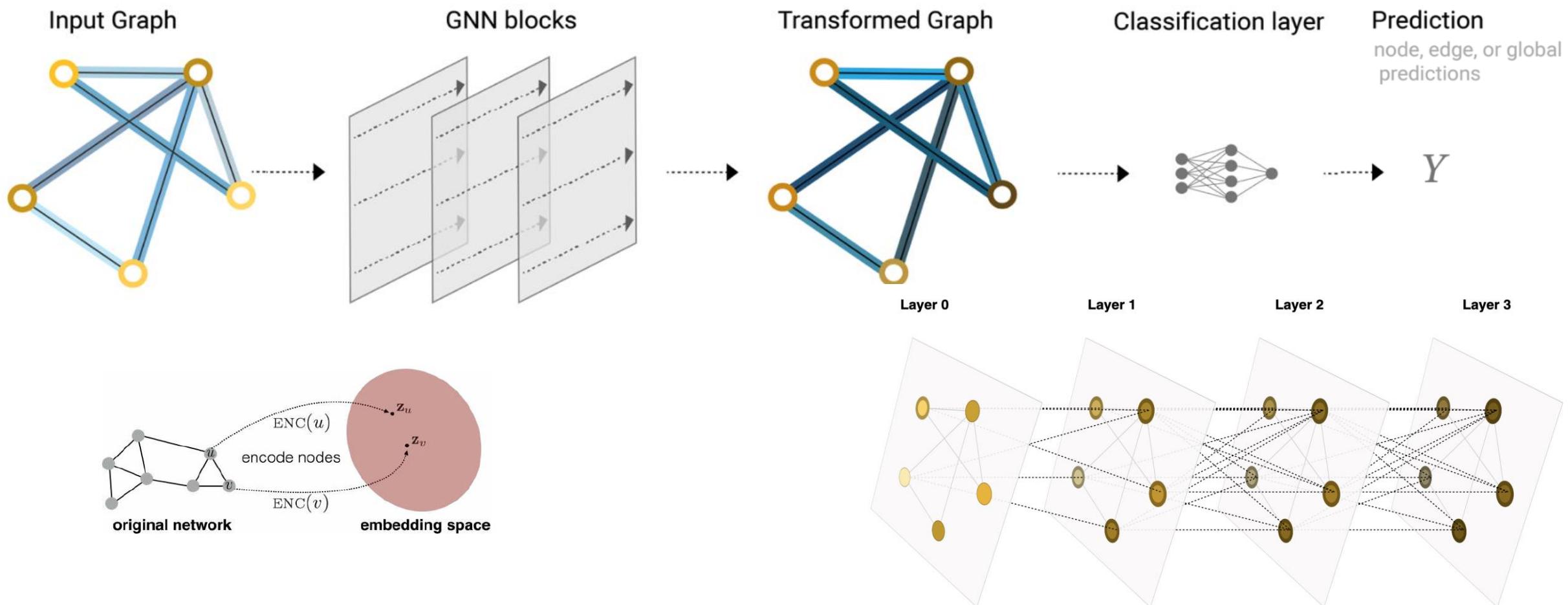


City A



City B

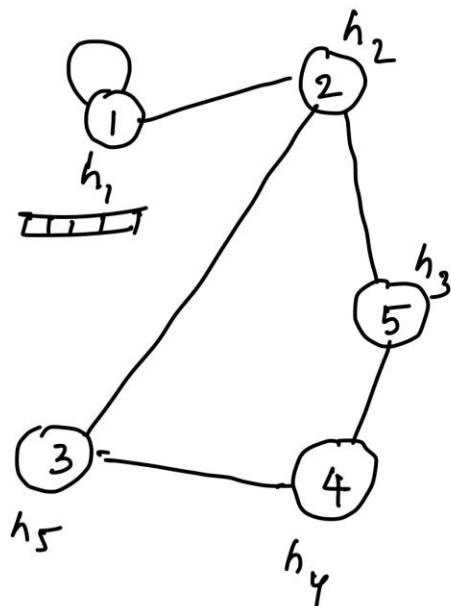
Graph neural networks (GNNs) are a family of neural networks that can operate naturally on graph-structured data.



$$h_i' = \sigma \left(\sum_{j \in \mathcal{N}(i)} w * h_j \right)$$

LET'S TAKE A DEEP DIVE INTO THE WORKING OF A SIMPLE GNN

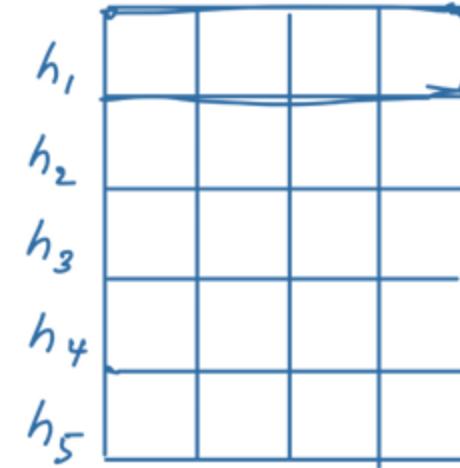
$$h_i' = \sigma \left(\sum_{j \in N(i)} W * h_j \right)$$



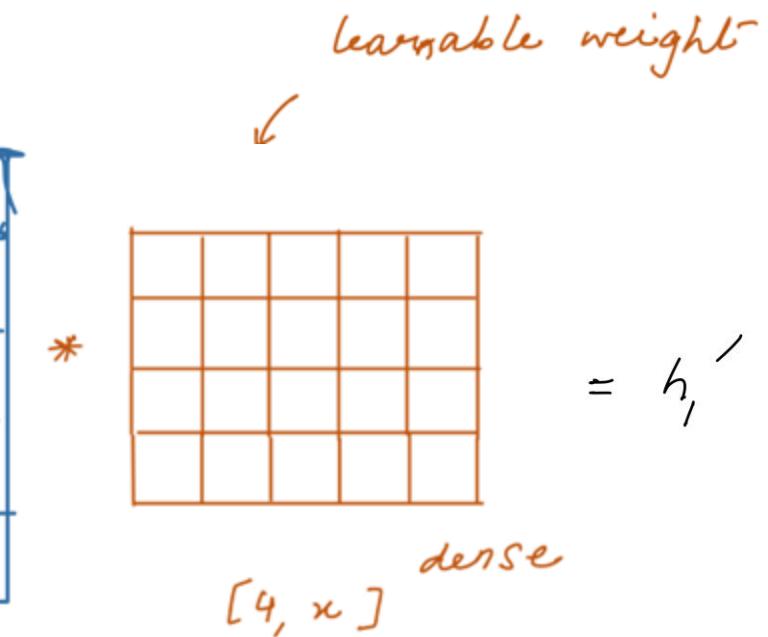
① $\begin{bmatrix} 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 \end{bmatrix}$

Adj Matrix

$[5, 5]$
(sparse)



$[5, 4]$



LET'S TAKE A DEEP DIVE INTO THE WORKING OF A SIMPLE GNN

$$h_i' = \sigma \left(\sum_{j \in \mathcal{N}(i)} W * h_j \right)$$



$$h_i^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}_i} \underbrace{\frac{1}{c_{ij}}}_{\text{normalizing term}} W^l * h_j^l \right)$$

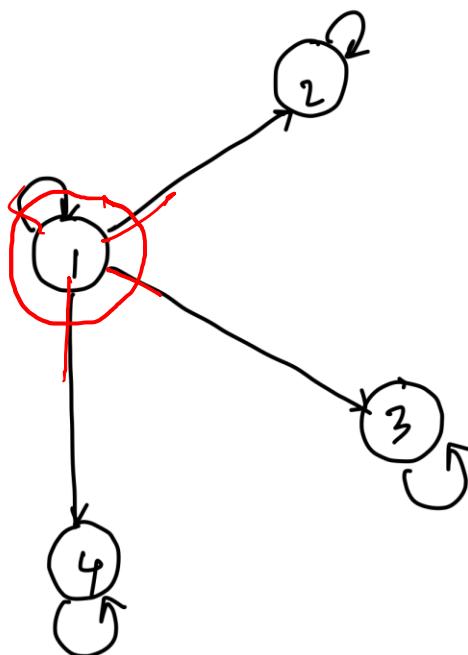
No normalizing term → After a few layers this will cause the node embeddings to explode.

Add a normalizing term based on the degree of the node.

GRAPH CONVOLUTIONAL NEURAL NETWORK

$$h_i^{l+1} = \sigma \left(\sum_{j \in N_i} c_{ij} W^l * h_j^l \right)$$

$$H^{(l+1)} = \sigma \left(\underbrace{\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}}_{\text{Norm}} H^{(l)} \underbrace{W^{(l)}}_{\text{projection}} \right)$$



$$[N, N] \xrightarrow{A} A = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & \ddots & & \\ 1 & & \ddots & \\ 1 & & & \ddots \end{bmatrix}$$

$$[N, N] \xrightarrow{A + I} \tilde{A} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix}$$

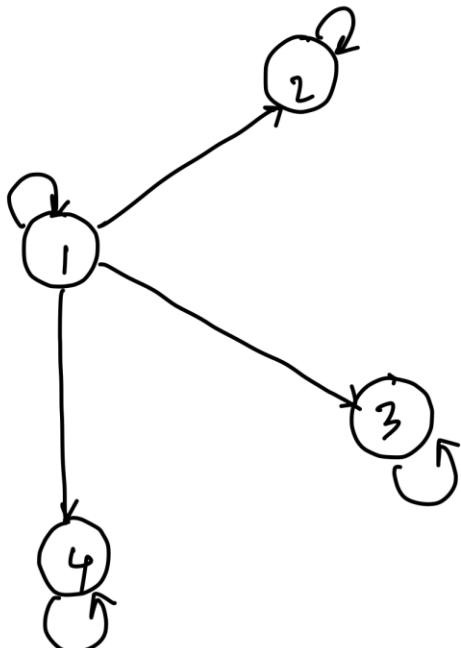
Degree

$$\xrightarrow{\tilde{D}} [N, N] = \begin{bmatrix} 4 & 0 & 0 & 0 \\ \cdot & 2 & \cdot & \cdot \\ 0 & \cdot & 2 & \cdot \\ \cdot & \cdot & \cdot & 2 \end{bmatrix}$$

GRAPH CONVOLUTIONAL NEURAL NETWORK

$$h_i^{l+1} = \sigma \left(\sum_{j \in N_i} c_{ij} \frac{1}{\|W^l * h_j^l\|} W^l * h_j^l \right)$$

$$H^{(l+1)} = \sigma \left(\underbrace{\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}}_{\text{Norm}} \underbrace{H^{(l)} W^{(l)}}_{\text{projection}} \right)$$



$$\tilde{A} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix}$$

$A + I$

$$\tilde{D} = \begin{bmatrix} 4 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 \end{bmatrix}$$

Degree

$$\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} \xrightarrow{\text{eq}} \frac{a_{ij}}{\sqrt{d_i d_j}} \xleftarrow{\text{Normalization}} \frac{a_{12}}{\sqrt{d_1 d_2}}$$

$$\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} = \begin{bmatrix} \frac{a_{11}}{\sqrt{d_1 d_1}} & \cdot & \cdot & \cdot \\ \cdot & \ddots & \cdot & \cdot \\ \cdot & \cdot & \ddots & \cdot \\ \cdot & \cdot & \cdot & \ddots \end{bmatrix}$$

LET'S TAKE A DEEP DIVE INTO THE WORKING OF A SIMPLE GNN

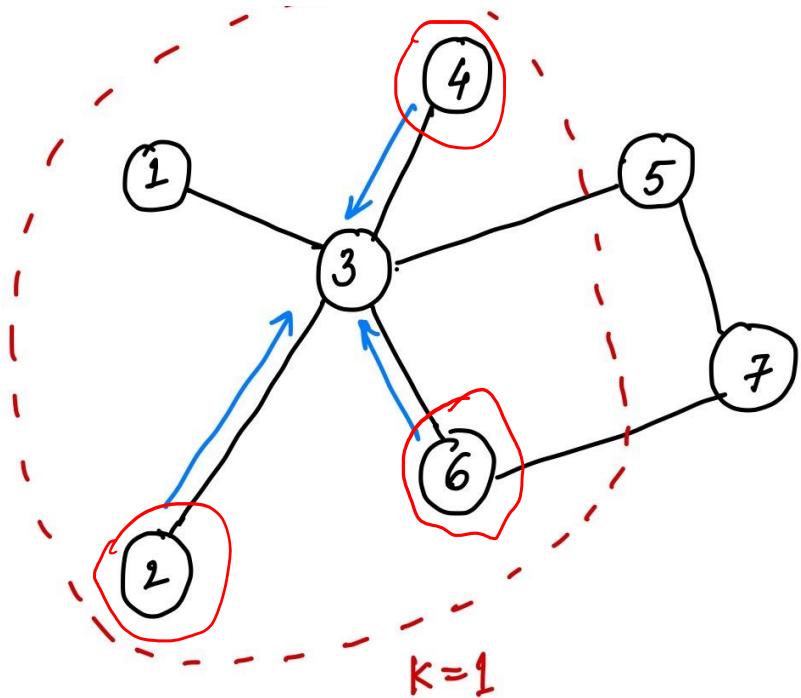
$$h_i' = \sigma \left(\sum_{j \in \mathcal{N}(i)} W * h_j \right)$$

$$h_i^{l+1} = \sigma \left(\sum_{j \in \mathcal{N}_i} c_{ij} \overset{\curvearrowleft}{W^l} * h_j^l \right) \rightarrow h_i^{l+1} = \sigma \left(W \cdot \left[\underset{j \in \mathcal{N}_i}{\text{AGGREGATOR}} \left(\{h_j^l\} \right), h_i^l \right] \right)$$

Update to a node embedding in GCN requires looking at all nodes. Very expensive operation

GraphSAGE: add sample and aggregator layer.

$$h_i^{l+1} = \sigma \left(W \cdot \left[\underset{j \in \mathcal{N}_i}{\text{AGGREGATOR}}(\{h_j^l\}), h_i^l \right] \right)$$



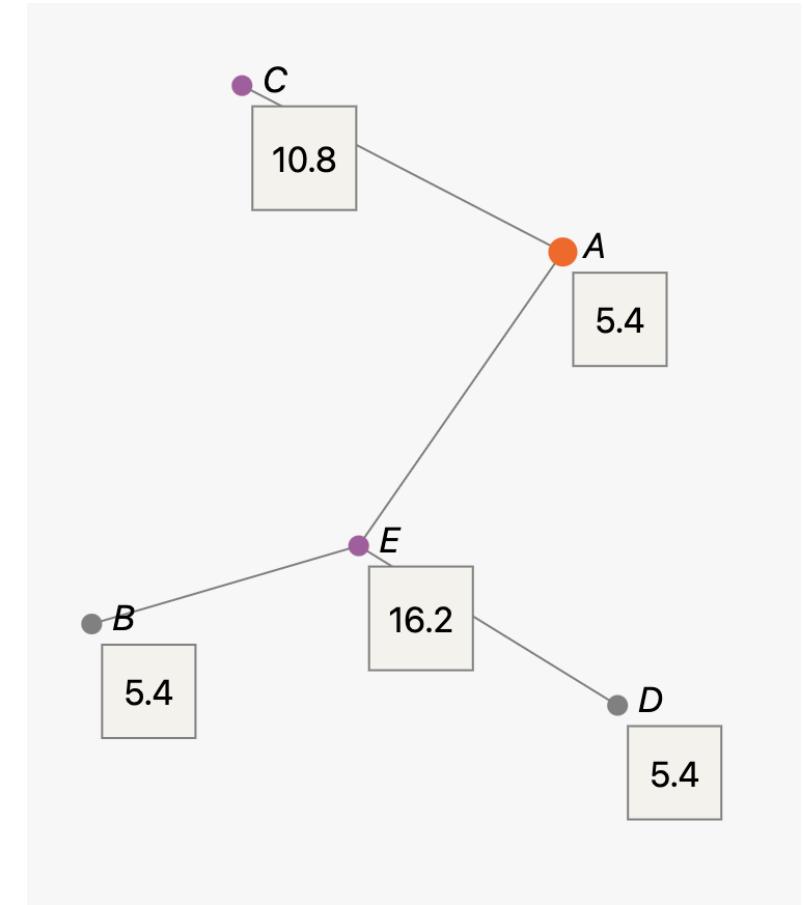
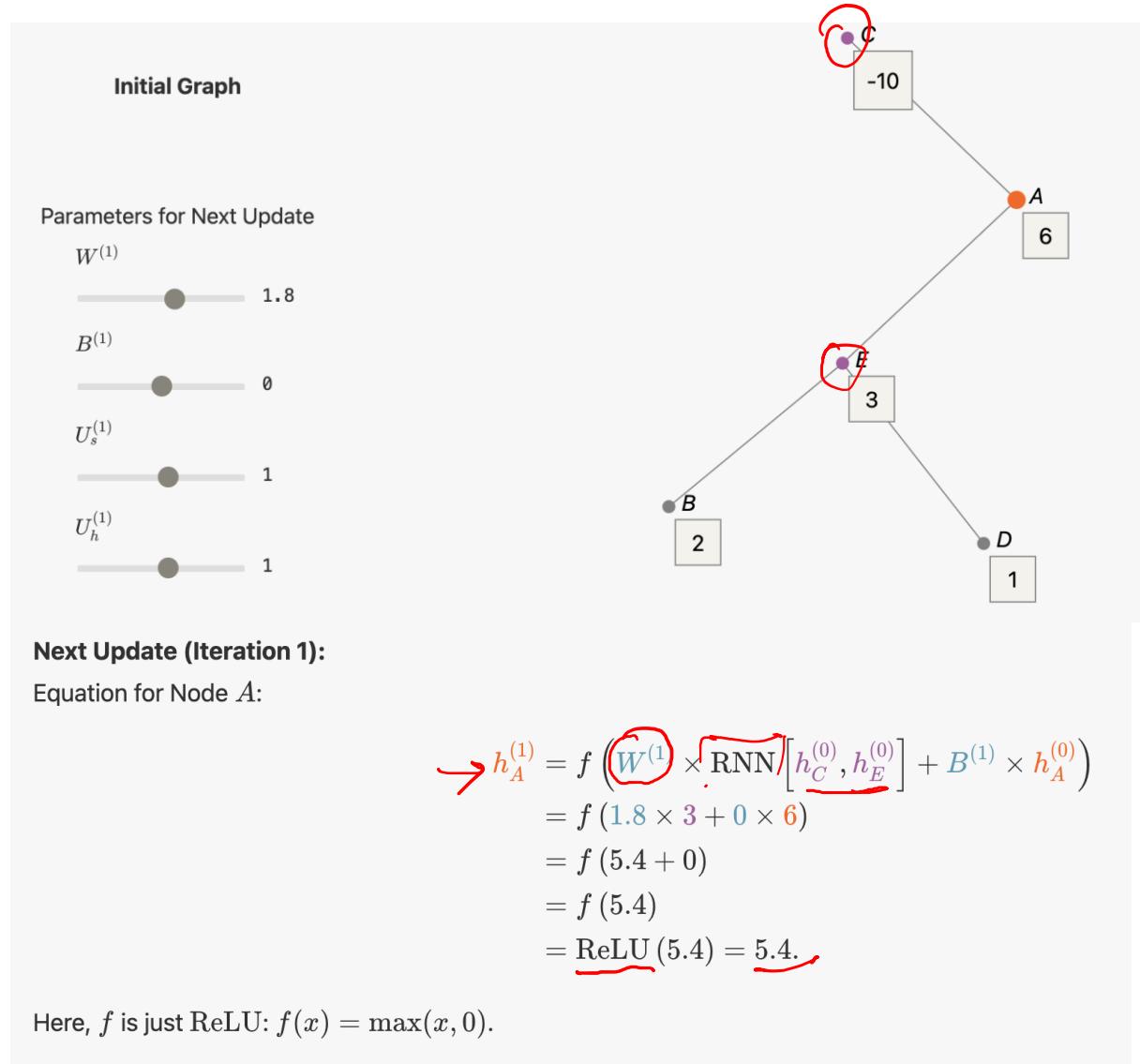
Types of aggregator

MEAN : $\underset{j \in \mathcal{N}_i}{\text{MEAN}}(\{h_j^l\}) \cup \{h_i^l\}$

POOLING : $\max(\{\sigma(W_{pool} h_j^l + b), \forall j \in \mathcal{N}_i\})$

And many more types of aggregators ...

GRAPH SAGE NETWORK



LET'S TAKE A DEEP DIVE INTO THE WORKING OF A SIMPLE GNN

$$h_i' = \sigma \left(\sum_{j \in N(i)} W * h_j \right)$$

$$h_i^{l+1} = \sigma \left(\sum_{j \in N_i} \underbrace{c_{ij}}_{\downarrow} W^l * h_j^l \right)$$

$$h_i^{l+1} = \sigma \left(W \cdot \left[\underset{j \in N_i}{\text{AGGREGATOR}} (\{h_j^l\}), h_i^l \right] \right)$$



$$\vec{h}_i' = \parallel_{k=1}^K \sigma \left(\sum_{j \in N_i} \alpha_{ij}^k W^k \vec{h}_j \right)$$

Random sampling does not provide **attention** to important nodes.

GAT: Add attention layer with learned parameter to weight the neighbors.

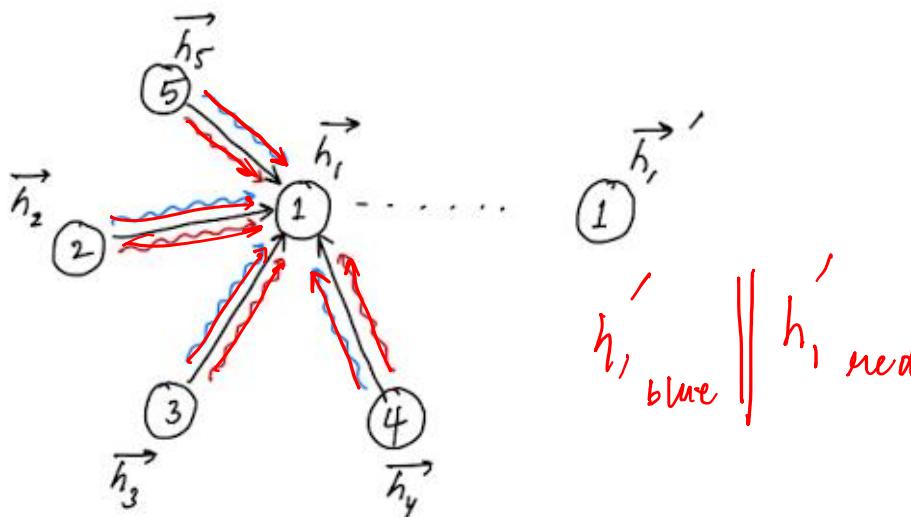
GRAPH ATTENTION NETWORK

Single head :

$$\vec{h}_i' = \sigma \left(\sum_{j \in \mathcal{N}_i} \alpha_{ij} \cdot W \vec{h}_j \right)$$

Multi-headed :

$$\vec{h}_i' = \parallel \sigma \left(\sum_{k=1}^K \sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \cdot W^k \vec{h}_j \right)$$



$$\alpha_{ij} = \text{softmax}(e_{ij})$$

↑
normalized attention coef.

→ Attention coefficients :

$$e_{ij} = a(W\vec{h}_i, W\vec{h}_j)$$

↑
attention mechanism

GRAPH ATTENTION NETWORK

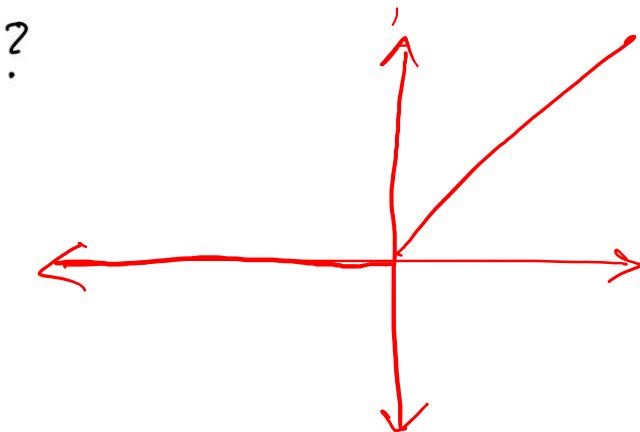
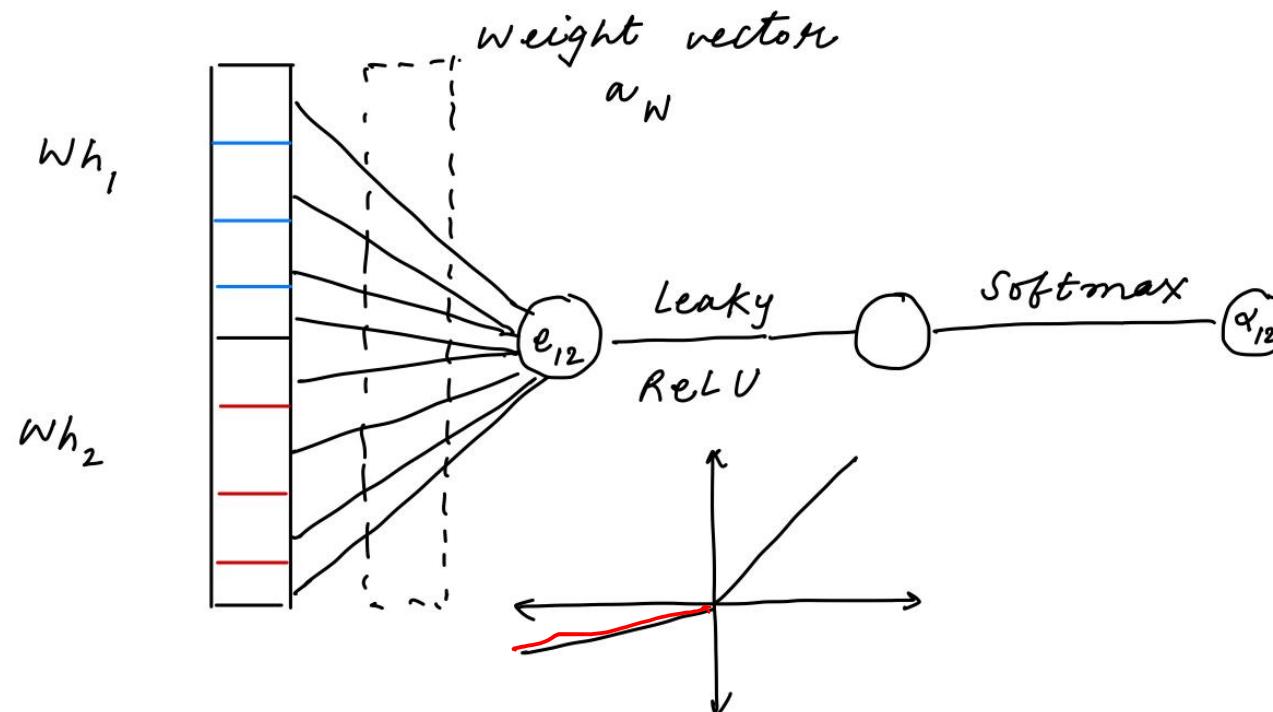
$$e_{ij} = \overset{\text{att. mechanism}}{a}(W h_i, W h_j)$$

attention coeff.

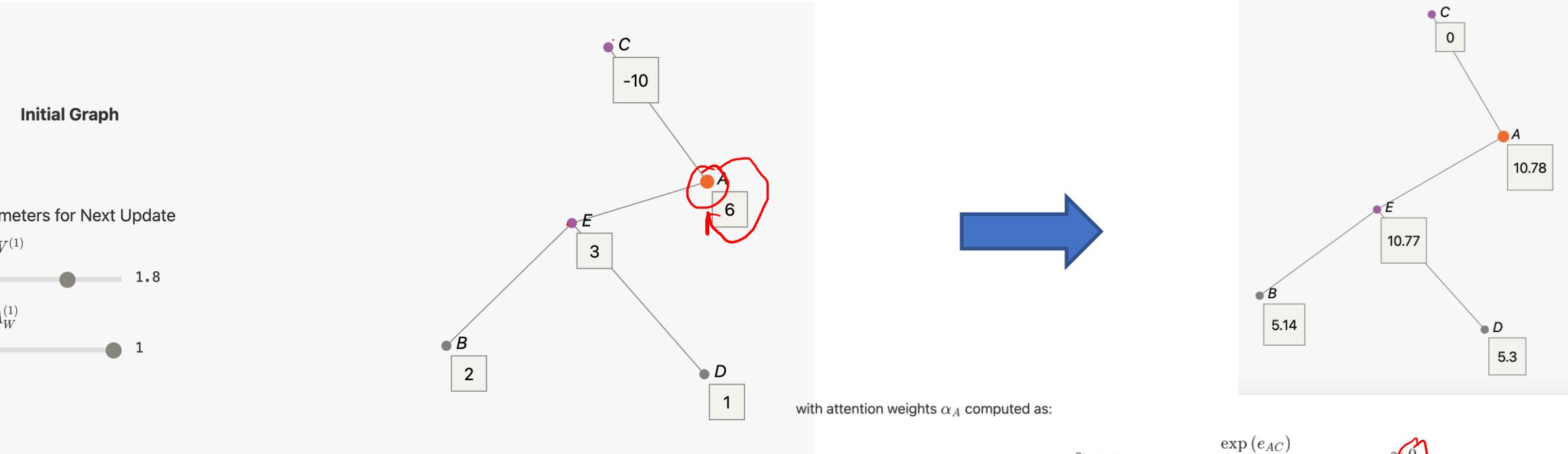
$$\alpha_{ij} = \text{softmax}(e_{ij})$$

normalized attention coeff.

What does attention mechanism 'a' do ?



GRAPH ATTENTION NETWORK



Next Update (Iteration 1):

Equation for Node A:

$$\begin{aligned}
 h_A^{(1)} &= f \left(W^{(1)} \times \left(\underbrace{\alpha_{AC} \times h_C^{(0)}}_{0 \times -10} + \underbrace{\alpha_{AE} \times h_E^{(0)}}_{0 \times 3} + \underbrace{\alpha_{AA} \times h_A^{(0)}}_{1 \times 6} \right) \right) \\
 &= f \left(W^{(1)} \times \left(\underbrace{0 \times h_C^{(0)}}_{0} + \underbrace{0 \times h_E^{(0)}}_{0} + \underbrace{1 \times h_A^{(0)}}_{6} \right) \right) \\
 &= f \left(W^{(1)} \times (0 \times -10 + 0 \times 3 + 1 \times 6) \right) \\
 &= f \left(W^{(1)} \times (0 + 0.01 + 5.97) \right) \\
 &= f (1.8 \times (5.99)) \\
 &= f (10.78) \\
 &= \text{ReLU}(10.78) = 10.78.
 \end{aligned}$$

$$\begin{aligned}
 \alpha_{AC} &= \frac{\exp(e_{AC})}{\exp(e_{AC}) + \exp(e_{AE}) + \exp(e_{AA})} \approx 0. \\
 \alpha_{AE} &= \frac{\exp(e_{AE})}{\exp(e_{AC}) + \exp(e_{AE}) + \exp(e_{AA})} \approx 0. \\
 \alpha_{AA} &= \frac{\exp(e_{AA})}{\exp(e_{AC}) + \exp(e_{AE}) + \exp(e_{AA})} \approx 1
 \end{aligned}$$

where the unnormalized attention weights e_A are given by:

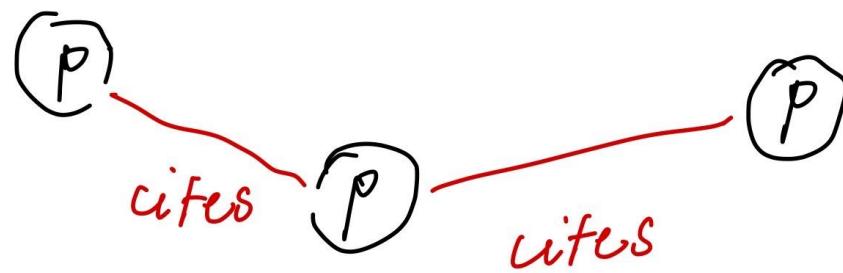
$$\begin{aligned}
 e_{AC} &= \text{ReLU} \left(A_W^{(1)} \left(W^{(1)} h_A^{(0)} + W^{(1)} h_C^{(0)} \right) \right) = \text{ReLU} (1 (1.8 \times 6 + 1.8 \times -10)) = \text{ReLU} (-7.2) = 0. \\
 e_{AE} &= \text{ReLU} \left(A_W^{(1)} \left(W^{(1)} h_A^{(0)} + W^{(1)} h_E^{(0)} \right) \right) = \text{ReLU} (1 (1.8 \times 6 + 1.8 \times 3)) = \text{ReLU} (16.2) = 16.2. \\
 e_{AA} &= \text{ReLU} \left(A_W^{(1)} \left(W^{(1)} h_A^{(0)} + W^{(1)} h_A^{(0)} \right) \right) = \text{ReLU} (1 (1.8 \times 6 + 1.8 \times 6)) = \text{ReLU} (21.6) = 21.6.
 \end{aligned}$$

We have omitted the superscripts on the attention weights for clarity.

Here, f is just ReLU: $f(x) = \max(x, 0)$.

<https://arxiv.org/pdf/1710.10903.pdf>

TYPES OF GRAPHS

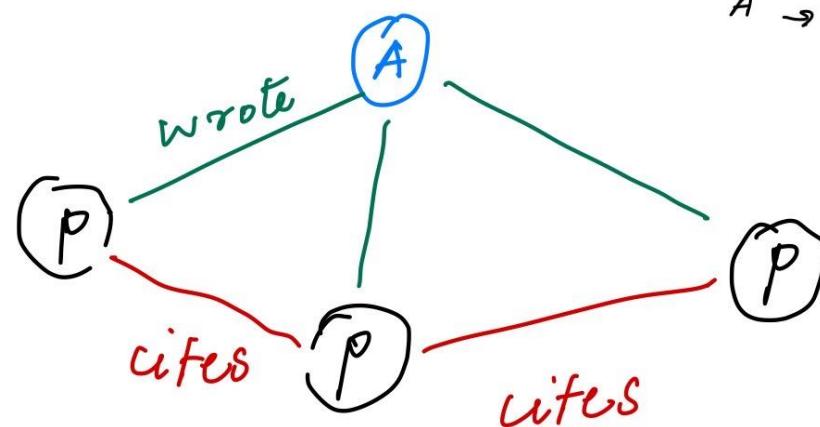


Apply GCN, GraphSAGE, GAT

(or many other models)

Homogeneous graph → 1 type of node, 1 type of relation

TYPES OF GRAPHS



$P \rightarrow$ Paper Node
 $A \rightarrow$ Author Node

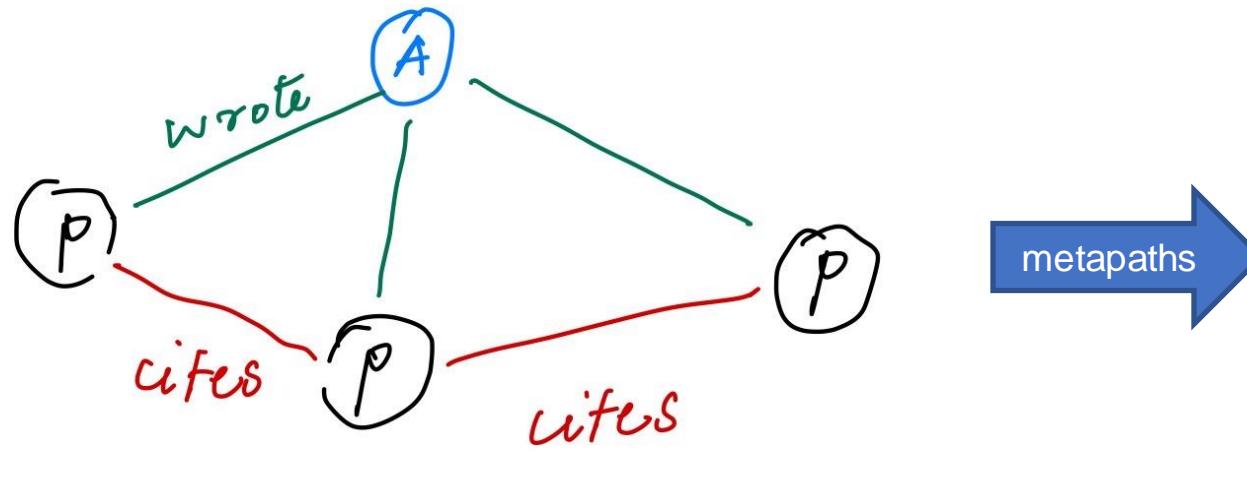
What models to apply ?

- ① Simplify
- ② Heterogeneous Graph Models
eg: HAN, HGAT

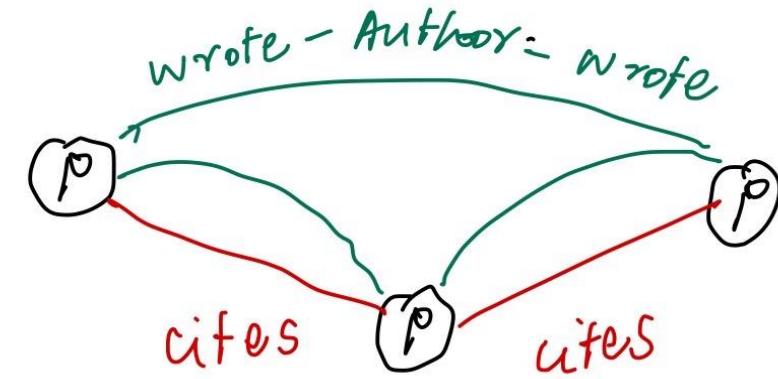
Heterogeneous graph

→ Multiple types of nodes,
Multiple types of relations.

TYPES OF GRAPHS



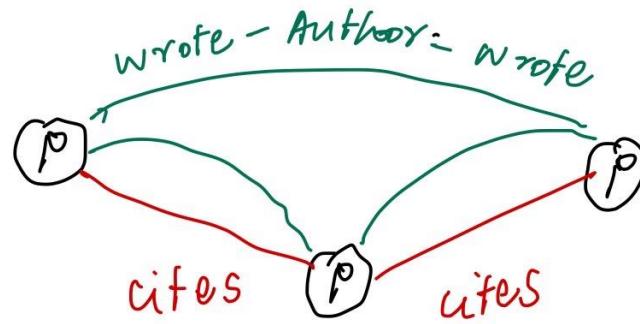
metapaths



Heterogeneous graph

Multirelational graphs

TYPES OF GRAPHS



Apply Relational GCN/GAT

Multirelational Graphs → 1 type of node

Multiple types of relations.

MOVING FROM HOMOGENEOUS MODELS TO RELATIONAL MODELS

$$h_i^{l+1} = \sigma \left(\sum_{j \in N(i)} W * h_j^l \right)$$

$$h_i^{l+1} = \sigma \left(\sum_{j \in N_i} \underbrace{\frac{1}{c_{ij}}}_{\text{norm}} W * h_j^l \right)$$

$$h_i^{l+1} = \sigma \left(W \cdot \left[\underset{j \in N_i}{\text{AGGREGATOR}}(\{h_j^l\}), h_i^l \right] \right)$$

$$\vec{h}_i^{l+1} = \left\| \sigma \left(\sum_{j \in N_i} \alpha_{ij}^k W \vec{h}_j^l \right) \right\|_k^k$$

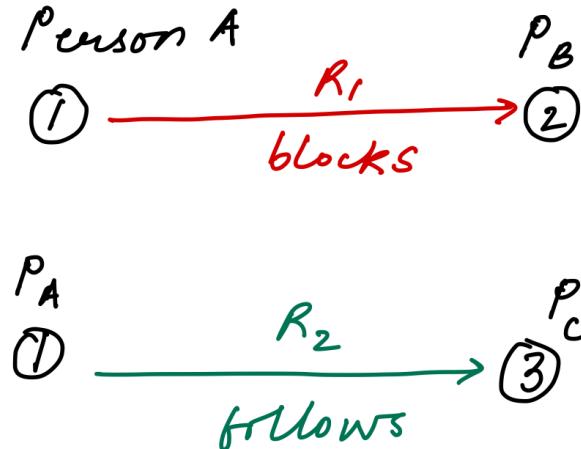
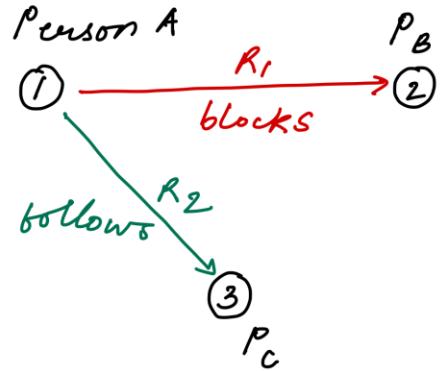
$$h_i^{l+1} = \sigma \left(\sum_{j \in N_i} \underbrace{\frac{1}{c_{ij}}}_{\text{norm}} W^l * h_j^l \right) \rightarrow h_i^{l+1} = \sigma \left(W_o^l h_i^l + \sum_{r \in R} \sum_{j \in N_i^r} \frac{1}{c_{ij}} W_r^l h_j^l \right)$$

The GCN model requires homogenous graph.
 GCNs use a single weight matrix for each layer
 which cannot utilize the information from
 different types of edges.

Heterogenous graph can be converted to
 relational graph with the meta path concept.

The R-GCN uses separate weight matrices for
 separate relations in order to uniquely use different
 relations to update the node embeddings.

RELATIONAL GRAPH CONVOLUTIONAL NEURAL NET



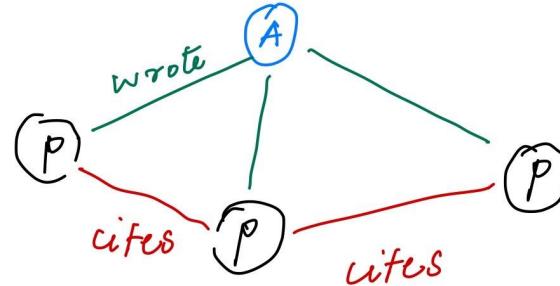
$$\sum_{j \in N_1} \frac{1}{c_{1j}} w_r h_j$$

$$\sum_{j \in N_i} \frac{1}{c_{ij}} w_r h_j$$

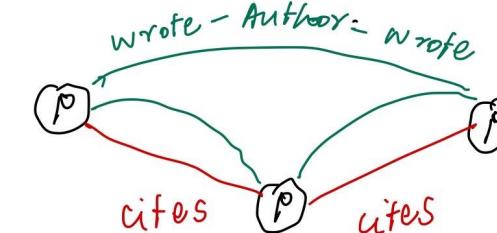
$$h_i' = \sigma \left(\sum_{r \in R} \sum_{j \in N_i} \frac{1}{c_{ij}} w_r h_j \right)$$

$$h_i' = \sigma \left(w_o h_i + \sum_{r \in R} \sum_{j \in N_i} \frac{1}{c_{ij}} w_r h_j \right)$$

RELATIONAL GRAPH CONVOLUTIONAL NEURAL NET



Heterogeneous graph



Multirelational graphs

$$h_i^{l+1} = \sigma \left(w_o^l h_i^l + \sum_{r \in R} \sum_{j \in N_i^r} \frac{1}{c_{ij}} w_r^l h_j^l \right)$$

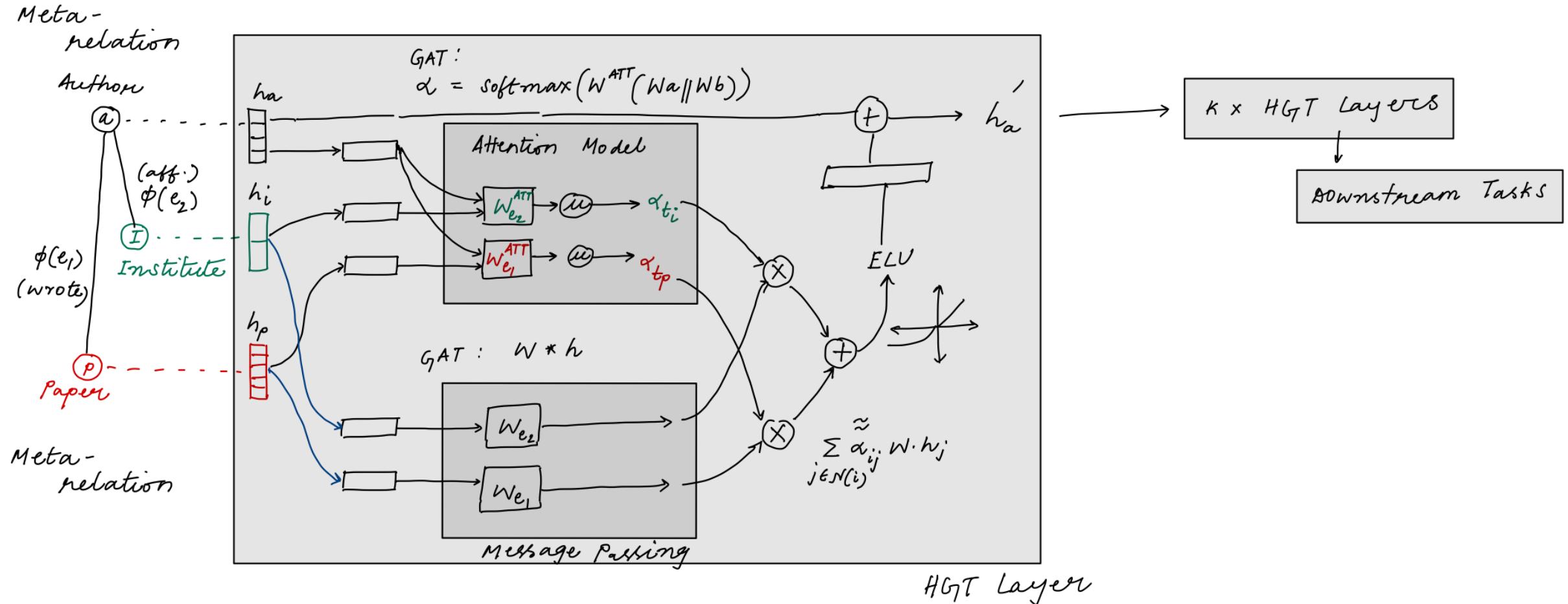


- > Metapath decomposition is done by GNN researcher who **may not have domain expertise** on the dataset
- > **Loss of information** due to simplification.

Heterogeneous Graph Transform

- > Introduce **meta-relations**
- > Use **GAT like computation**.

HETEROGENOUS GRAPH TRANSFORMER



LOSS FUNCTION DEPENDING ON TASK

- **Node Classification:**

$$\mathcal{L}(y_v, \hat{y}_v) = - \sum_c y_{vc} \log \hat{y}_{vc}.$$

where \hat{y}_{vc} is the predicted probability that node v is in class c .

So, $-\sum_c y_{vc} \log \hat{y}_{vc}$ is the total loss of all the classes of a multiclass classification.

e.g.:

	TARGET	PREDICTION	$\mathcal{L}(y_{vc}, \hat{y}_{vc})$
Science	$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$	$\rightarrow 0.356$
Math	$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$	$\rightarrow 0$
English	$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$	$\begin{bmatrix} 0.7 \\ 0.2 \\ 0.1 \end{bmatrix}$	$\rightarrow 0$

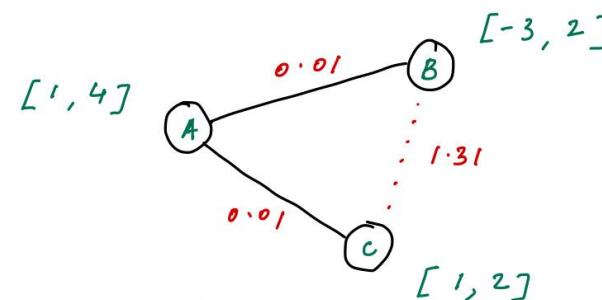
$$\mathcal{L}(y_v, \hat{y}_v) = 0.356$$

LOSS FUNCTION DEPENDING ON TASK

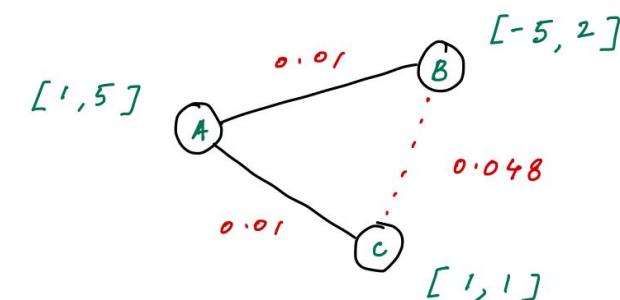
- Link Prediction: $\mathcal{L}(y_v, y_u, e_{vu}) = -e_{vu} \log(p_{vu}) - (1 - e_{vu}) \log(1 - p_{vu})$
$$p_{vu} = \sigma(y_v^T y_u)$$

where e_{uv} is 1 when there is an edge between u and v node.

Let's see an example



$$\begin{aligned} P_{AB} &= 0.99 \\ P_{AC} &= 0.99 \\ P_{CB} &= 0.73 \end{aligned}$$



$$\begin{aligned} P_{AB} &= 0.99 \\ P_{AC} &= 0.99 \\ P_{CB} &= 0.047 \end{aligned}$$

CHALLENGES TO GRAPH NEURAL NET COMMUNITY

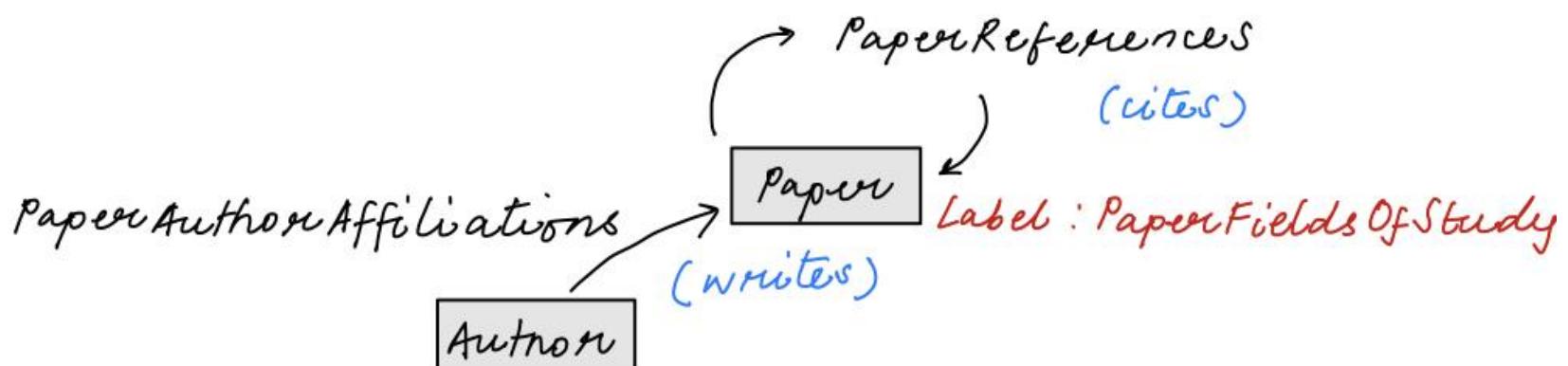
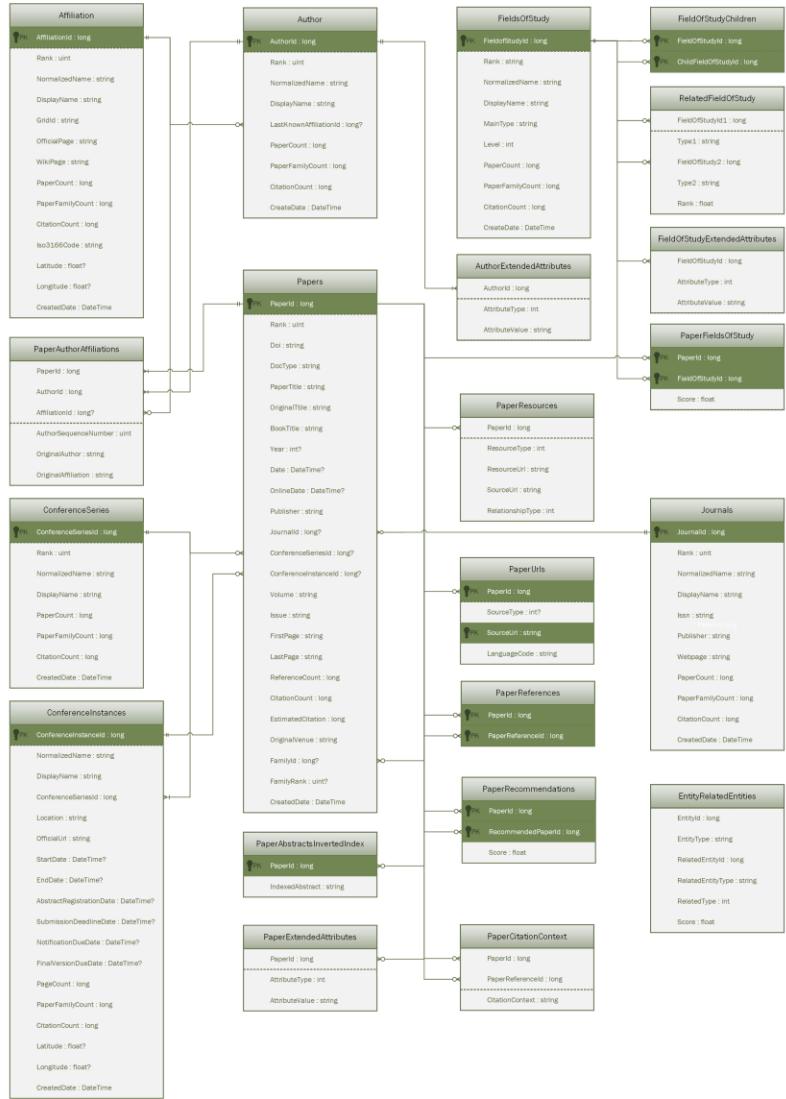


Dataset	Date	Node emb size	#nodes (millions)	#edges (millions)	#labeled nodes (millions)
OBGN-papers	2020	128*	111	1,615	1.4
MAG240M	2021	768	260	1,300	1.4
PinSAGE dataset	2018	128-2K (avg. 1K)	3000	18,000	UNK

- Large gap between publicly available graph datasets and ones used by industry.
- Our goal is to propose a new dataset that will help both system designers and GNN researchers in two ways:
- Given a dataset schema, propose a methodology to generate arbitrary sized graphs (homogenous or heterogenous) and node embeddings with prescribed number of nodes, edges and relations.
- Provide a dataset with synthetic node embeddings for system developers and another dataset with node embeddings generated using NLP methods for GNN researchers and system developers.

* → paper does not provide information. Calculated via dataset size.

ILLINOIS GRAPH BENCHMARK (IGB)

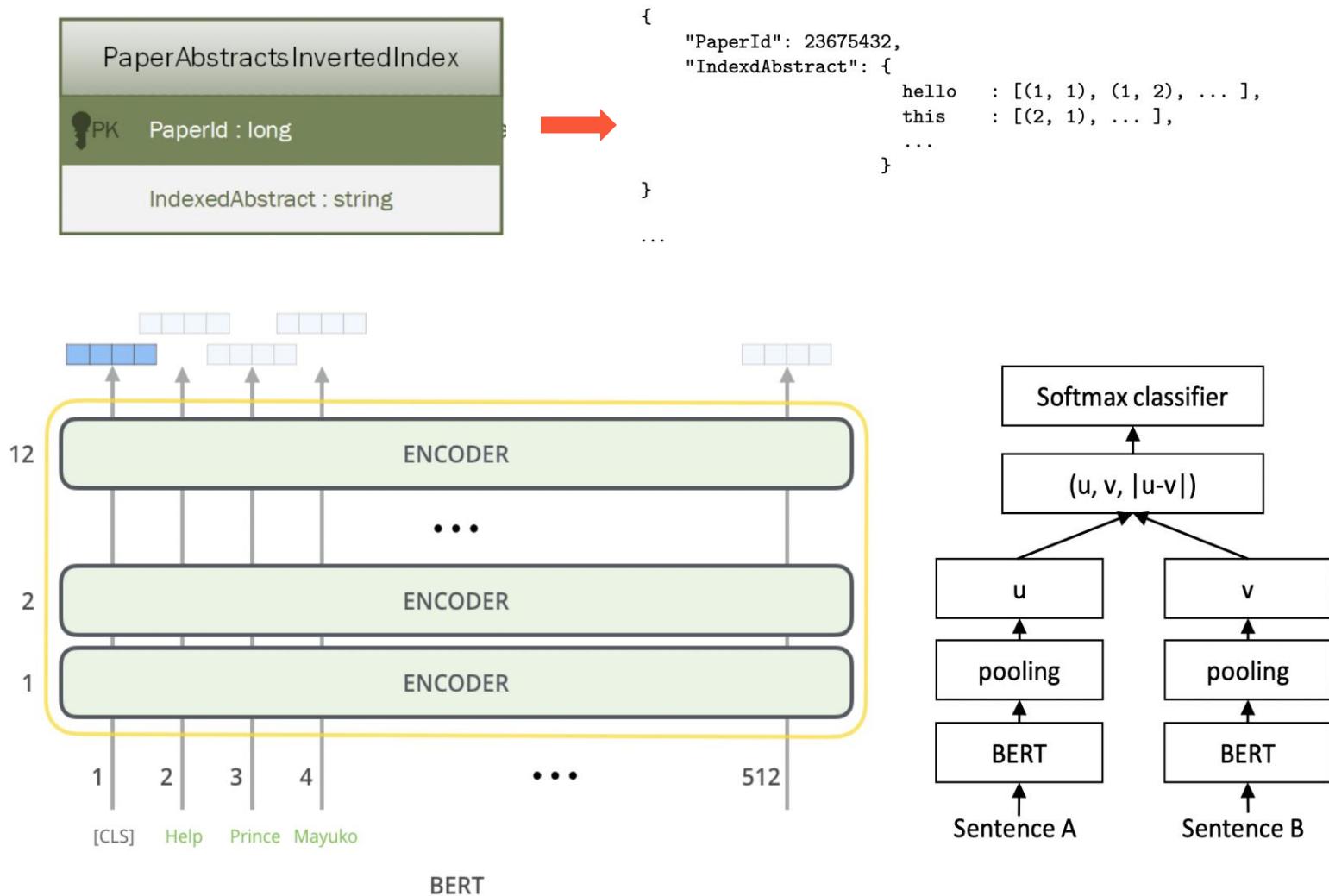


GENERATION OF NODE EMBEDDINGS

From the MAG dataset we extracted abstracts for each paper which is stored in an inverted index format.

After converting it into text and cleaning the data, we used the Sentence-BERT model to generate node embeddings for each paper node.

The BERT model takes the entire abstract and tokenizes each word using a layer of Encoders which is then pooled together to get a single embedding using the Sentence-BERT model.



Alammar, J (2018). The Illustrated Transformer [Blog post]. Retrieved from <https://jalammar.github.io/illustrated-transformer/>

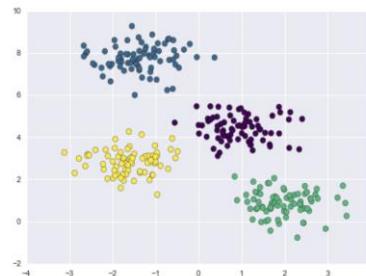
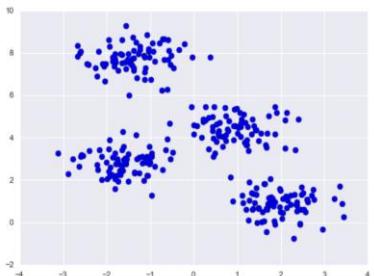
HOW WE CREATED THE LABELS

Challenges: The largest graph dataset that is available to the public today has **only 1.4 million labeled nodes**.

1. To over come this limitation in addition to the arxiv labels, we have incorporated an additional database namely the Semantic Scholar to get an **additional 100 million labeled nodes** out of 260 million nodes.

Journals	
PK	JournalId : long
Rank : unit	
Normalized Name : string	
Display Name : string	

ConferenceInstances	
PK	ConferenceInstanceId : long
Normalized Name : string	
Display Name : string	
Conference Series Id : long	
Location : string	



```
{  
  "id": "4cd223df721b722b1c40689caa52932a41fcc223",  
  "fieldsOfStudy": ["Computer Science"],  
}
```

2. To generate even more labels for large prediction models we used an unsupervised K-Means algorithm using the journal and conferences information of each paper to cluster them into K classes. This works in two-folds: it allows the user to train on **more data** and allows the user to select the number of classes to get a finer or coarser classification.

CURRENT STATUS

Dataset	Date	Type	Node emb size	#nodes (millions)	#edges (millions)	#labeled nodes (millions)
OBGN-papers	2020	Real	128*	111	1,615	1.4
MAG240M	2021	Real	768	260	1,300	1.4
IGB v1	2022	Real/Synthetic	128 - 4kB	267	1,900	~130+
IGB v2	2022	Real/Synthetic	variable	~600+	~ 3,000+	~200+
PinSAGE dataset	2018	Real	128-2K	3000	18,000	NA

- Uses Sentence-BERT model to generate “real” node embeddings .
- Synthetic embedding vs real embedding:
 - Fundamentally GNN's find structural information of the graph to improve the embeddings and have no idea about the node's information.
 - Synthetic graph dataset would be useful to test computation and optimization the results
 - Synthetic graph datasets do not have any real-world significance for downstream tasks. For example, GAT would be completely useless for a synthetic dataset
- Then we will run GNN models to establish baseline results.

- Overall, our goal is to help the growth of the Graph Neural Net community
 - By giving GNN researchers access to larger labelled real graph datasets to **test and develop more robust models**.
 - By providing system researchers access to large synthetic datasets to **improve the hardware and systems** so we can process large amount of data more efficiently.

We are continuously running baseline experiments to confirm the working on the dataset and prove that larger datasets lead to higher baseline accuracy in supervised tasks.

If you have any questions about Graph Neural Nets or the Illinois Graph Benchmark at a later point, feel free to reach out to me at akhatua2@illinois.edu

Thank you!

Please feel free to ask me any questions



The Grainger College of Engineering

UNIVERSITY OF ILLINOIS URBANA-CHAMPAIGN